

Collaboration Policy: You are encouraged to collaborate with up to 4 other students, but all work submitted must be your own *independently* written solution. List the computing ids of all of your collaborators in the `collabs` command at the top of the tex file. Do not share written notes, documents (including Google docs, Overleaf docs, discussion notes, PDFs), or code. Do not seek published or online solutions for any assignments. If you use any published or online resources (which may not include solutions) when completing this assignment, be sure to cite them. Do not submit a solution that you are unable to explain orally to a member of the course staff. Any solutions that share similar text/code will be considered in breach of this policy. Please refer to the syllabus for a complete description of the collaboration policy.

Collaborators: Yining Liu(yl2nr), Yuhao Niu(bhb9ba)

Sources:

PROBLEM 1 Solving Recurrences

Prove a (as tight as possible) O (big-Oh) asymptotic bound on the following recurrences. You may use any base cases you'd like.

1. For the following two recurrences, it may be helpful to draw out the tree. However, you should prove the asymptotic bound using induction.

- $T(n) = T(\frac{n}{2}) + T(\frac{n}{4}) + T(\frac{n}{8}) + n$

Solution:

Proof.

We want to show that the function $T(n) = T(\frac{n}{2}) + T(\frac{n}{4}) + T(\frac{n}{8}) + n \in O(n \log n)$.

For the base case, consider $n = 1$, $T(1) = c$, and it is a constant.

Next, we assume that for some positive integer k , $T(n) \in O(n \log n)$ for all n up to 2^k .

Then, we want to prove that $T(n) \in O(n \log n)$ for $n = 2^{k+1}$.

Using recurrence relation:

$$T(2^{k+1}) = T(2^k) + T(\frac{2^k}{2}) + T(\frac{2^k}{4}) + 2^{k+1}$$

Based on the induction hypothesis: $T(2^k) \in O(2^k \log 2^k) = O(k \cdot 2^k)$

Thus, $T(2^{k-1}) \in O((k-1) \cdot 2^{k-1})$, and $T(2^{k-2}) \in O((k-2) \cdot 2^{k-2})$

So, we can write that for some $m > 0$:

$$T(2^{k+1}) \leq m(k \cdot 2^k + (k-1) \cdot 2^{k-1} + (k-2) \cdot 2^{k-2})$$

$$T(2^{k+1}) \leq m \cdot 2^{(k+1)}(7k - 6)$$

Choose $c = 7m$, we have:

$$T(2^{k+1}) \leq 7m \cdot (2^{k+1} \cdot (k+1))$$

Because $7m$ is a constant multiplier, we can conclude that:

$$T(2^{k+1}) \in O(2^{k+1} \log 2^{k+1})$$

Hence we show that $T(n) \in O(n \log n)$ for all positive integers n . \square

- $T(n) = 2T(\frac{n}{3}) + T(\frac{n}{6}) + n$

Solution:

Proof.

We want to show that the function $T(n) = 2T(\frac{n}{3}) + T(\frac{n}{6}) + n \in O(n \log n)$.

For the base case that $n = 1, T(n) = c$, which is a constant.

Now we assume for some positive integer $k, T(n) \in O(n \log n)$ for all n up to 3^k .

So we want to prove that $T(n) \in O(n \log n)$ for $n = 3^{k+1}$.

Using recurrence relation:

$$T(3^{k+1}) = 2T(3^k) + T(\frac{3^k}{2}) + 3^{k+1}$$

Based on the induction hypothesis: $T(3^k) \in O(3^k \log 3^k) = O(k \cdot 3^k)$

So, we can write that for some $m > 0$:

$$T(3^{k+1}) \leq m \cdot (3k \log 3^k + \frac{3^k}{2} \log \frac{3^k}{2} + 3^{k+1})$$

$$T(3^{k+1}) \leq m \cdot 3^k \cdot (\log 3^k + \frac{1}{2} \log \frac{3^k}{2} + 3)$$

Choose $c = 3m$, we have:

$$T(3^{k+1}) \leq 3m \cdot (3^{k+1} \cdot \log 3^{k+1})$$

Because $3m$ is a constant multiplier, we can conclude that:

$$T(3^{k+1}) \in O(3^{k+1} \log 3^{k+1})$$

Hence we show that $T(n) \in O(n \log n)$ for all positive integers n . \square

2. For the following recurrence relations, indicate: (i) which case of the Master Theorem applies (if any); (ii) justification for why that case applies (if one does) i.e., what is $a, f(n), \epsilon$, etc; (iii) the asymptotic growth of the recurrence (if any case applies).

- $T(n) = 7T(\frac{n}{5}) + n \log n$

Solution:

(i)

Case 1 applies.

(ii)

$$a = 7, b = 5, \log_5 7 > 1$$

$$f(n) = n \log n, f(n) \in O(n^{(\log_5 7) - \epsilon}), \text{ for any } \epsilon < (\log_5 7 - 1).$$

(iii)

$$T(n) \in \Theta(n^{\log_5 7})$$

- $T(n) = 3T(\frac{n}{3}) + n \log n$

Solution:

(i)

No case applies.

(ii)

$$a = 3, b = 3, \log_3 3 = 1$$

$$f(n) = n \log n, \forall \epsilon > 0, f(n) \notin O(n^{\delta - \epsilon}), f(n) \notin \Omega(n^{\delta + \epsilon}), \text{ and } f(n) \notin \Theta(n^{\delta})$$

PROBLEM 2 Climate History

Scientists call *paleoclimatologists* study the history of climate on earth before instruments were invented to measure temperatures, precipitation, etc. They try to reconstruct climate history using data found from analyzing rocks, sediments, tree rings, fossils, ice sheets, etc.

A group is trying to better understand periods of low precipitation (e.g. dryness) for a time period that spans many thousands of years. Using various data sources, they have assigned a value d_i for the relative dryness for each time-unit (let's say it's measured for every century) for this very long time period. For each time-unit i , one of five values has been assigned as d_i :

- -3 for very high precipitation
- -1 for high precipitation
- 0 for normal (or unknown)
- 1 for low precipitation
- 3 for very low precipitation

The scientists want to find the score for the driest period of consecutive years in their data. Because the effects of a drought are cumulative, the score for a period starting at time i and ending at time j will be the value when all scores from i and j are added. Consider this example with $n = 16$:

d_i	0	3	3	0	-3	-3	1	0	3	-1	-1	0	1	1	3	0
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

The period [1:2] looks high with a score of 6. But the highest score is 7 for this period [6:14]. Note there are some negative values in the middle of this period, but the high values around those make it worth having a period that includes those negative values. (The period [6:15] also has score 7. There can be more than one period with the largest score, but you just need to return the highest score, so the existence of multiple periods with the same best score doesn't matter.)

There are many ways to solve this problem, including a $\Theta(n^2)$ brute-force approach that calculates scores for all combinations of starting and ending values for a period and keeps track of the largest score. You need to do better than that. (Though if you were ever to code your solution, coding the brute-force approach is a good way to test your algorithm. Coding is **not** required for this problem.)

What you need to do: Describe an algorithm (clearly in words or in pseudocode) that uses a divide and conquer algorithm that solves this problem in $\Theta(n \log n)$ time. Your description should make it clear what the base case is, what work is done in dividing and in combining, and what recursive subproblems are solved. Along with your algorithm's description, give a brief but convincing argument that the time complexity is $\Theta(n \log n)$.

The inputs to your algorithm will be n , the number of samples, and a list d that stores n values, where each is either -3, -1, 0, 1 or 3. Your algorithm will output the largest dryness score for a period in the data (as described above).

Solution:

We use divide and conquer to look for the highest dryness score over the period. We will recursively divide the array into smaller arrays till that each subarray has only one day and its score of dryness. Then we calculate the score and combine subarrays together while we also combine(add) the dryness score and keep track of the highest dryness score.

Divide:

1. Split the input array into two part, *Left* and *Right*.
2. Recursively divide the array into two halves until each subarray has only one value (single day). Record the midpoint while split.

Conquer:

1. After all the subarrays contain only one value (it is the base case), return the value of the dryness score.
2. For all the subarrays, we record the max dryness score in that array as *MaxDry*. For the base case, *MaxDry* = element in the array (because the base case has only one element, one dryness score).
3. For other cases (non-base case), *MaxDry* = $\max(\text{MaxDry of the left half}, \text{MaxDry of the right half}, \text{MaxDry of combined})$. It will be calculated during the Combine part.

Combine:

1. To calculate *MaxDryLeft* and *MaxDryRight*, we do this when we combine each subarrays. When subarrays are combined, they compare the current *MaxDry* to the sum of two *MaxDry* from two subarrays, and we store the larger one among them. When we combine all subarrays till we have only two subarrays *Left* and *Right*, the dryness score is then *MaxDryLeft* and *MaxDryRight*.
2. Then we combine *MaxDryLeft* and *MaxDryRight* to form a *MaxDryCombined*. We do this when we combine two parts of subarrays into one longer array (when we combine *Left* and *Right* into one piece).
4. Then, we have three values, *MaxDryLeft*, *MaxDryRight*, and *MaxDryCombined*. We return the *FinalDryness* as the max value among them.
5. The final combination back to the original array will produce the maximum dryness score of any possible period since we record the max value throughout the process.

The time complexity is $O(n \log n)$ because by dividing, the recursion tree has a depth of $\log n$. And each conquer and combine step is a linear process, so is at most n . Thus the final time complexity is $O(n \log n)$.

PROBLEM 3 *Fast Transformations*

Computer graphics software typically represents points in n dimensions as $(n + 1)$ -dimensional vectors. To make transformations on the points (e.g., rotating a modelled figure, zooming in, or making the figure appear as though seen through a fish-eye lens), we use a $(n + 1) \times (n + 1)$ matrix T which defines the transformation, and then we multiply each vector by this matrix to transform that point. That is, for vector v , the transformed vector is $v' = T \times v$.

Let's say we are developing software for very high dimension graphics (n dimensions), and we have a transformation T that we would like to apply to a particular point n times. Develop an algorithm which can multiply this $(n + 1)$ -dimensional point by T (the $(n + 1) \times (n + 1)$ transformation matrix) n times in $o(n^3)$ (little-oh of n^3) time. Prove this run time.

Solution:

We apply divide and conquer so that we put the total multiplication into smaller pieces (in half) of multiplication and recursively calculate the matrix until the base case. After we get the base case, we multiply vector v .

```
def Transformation(T,v,n)
    if n == 1:
        return T * v
    if n % 2 == 0:
        return Transformation(T*T,v,n/2)
    else:
        return T * Transformation(T*T,v,n/2)
```

The calculation of the product of matrices is strictly less than n^3 when we apply Strassen's algorithm to calculate it. We only have $\log n$ recursive cases, so the final time complexity is $o(n^3)$.

PROBLEM 4 *Receding Airlines*

You have been hired to plan the flights for Professor Floryan's brand new passenger air company, "Receding Airlines." Your objective is to provide service to n major cities within North America. The catch is that this airline will only fly you East.

You recognize that in order to enable all your passengers to travel from any city to any other city (to the East) with a single flight requires $\Omega(n^2)$ different routes. Prof. Floryan says that the airline cannot be profitable when supporting so many routes. Another option would be to order the cities in a list (from West to East), and have flights that go from the city at index i , to the city at index $i + 1$. This requires $\Theta(n)$ routes, but would mean that some passengers would require $\Omega(n)$ connections to get to their destination.

1. Devise a compromise set of routes which requires no passenger have more than a single connection (i.e. must take at most two flights), and requires no more than $O(n \log n)$ routes. Prove that your set of routes satisfies these requirements.
2. After a few years, passengers start demanding routes from East to West, and you decide to support new routes from East to West (in addition to supporting routes from West to East). Show that with routes in *both* directions, it is possible to connect all n cities with just $O(n)$ routes such that no passenger needs more than a single connection to get to their destination.

Note: In class so far, we've used recurrence relations to count an algorithm's time complexity, i.e., how many basic operations are executed. Here we're using one to count the number of flights. While this is a bit different than measuring an algorithm's time-complexity, we can still use a divide and conquer approach and a recurrence relation to answer these questions.

Solution:

For 1:

We can build a binary tree that has a middle airport as the root. Every airport on the left side of the tree (root) is on the west side of the middle airport. Then, every node (airport) on the right side of the tree (middle airport) is on the east side of the middle airport.

So all the branches of the binary tree will be the flying route for the company. We will also add routes from the root to every leaf of the tree. So all airports are connected to their root airport. Hence, airports at one level can reach every airport at the lower level.

Assume we have n airports. Then, from the middle, airports will have $n - 1$ routes. We then will apply divide and conquer to recursively calculate the number of routes at each lower level until we count every airport (node).

$$T(n) = 2T\left(\frac{n}{2}\right) + (n - 1)$$

And for base case:

$$T(1) = 0$$

To prove that it meets the requirement that flying from any airport to another takes at most one transition, we will take any two A_i and A_j airports. We just need to find their common ancestor A_a in the tree. So the airport on the west can fly to A_a and then from A_a to the airport on the east.

To prove the time complexity is $O(n \log n)$, we use Master Theorem:

$$a = 2, b = 2, \log_2 2 = 1, f(n) = n - 1 \in \Theta(n^1) = \Theta(n)$$

This is then the Case 2, leading to the implication that:

$$T(n) \in \Theta(n \log n) \Rightarrow T(n) \in O(n \log n)$$

$T(n)$ is the total number of flights, so the total number of flights is $O(n \log n)$.

For 2:

Now there is demand to go from west to east. When we want to ensure no one will have more than one transition, we need a connection airport to connect with all other airports. It is the middle airport in Case 1 since it is linked to all other airports. Now we need to add extra routes for each route to make it two directions. So any passenger can fly to the connection and then go to their destination. Now since we need a two-way route, we take the route as 2. Using the recursive steps to calculate, we get:

$$T(n) = 2T\left(\frac{n}{2}\right) + 2$$

To prove that flying from one airport to another needs at most 1 transaction, 2 flights, and the total number of routes is $O(n)$, we use Master Theorem:

$$a = 2, b = 2, \log_2 2 = 1, f(n) = 2 \in O(n^{1-0.01}) = O(n^{0.09})$$

So it is Case 1:

$$T(n) \in \Theta(n) \Rightarrow T(n) \in O(n)$$

$T(n)$ is the total number of flights, so the total number of flights is $O(n)$.

PROBLEM 5 *Bazinga!*

Theoretical Physicist Sheldon Cooper has decided to give up on String Theory in favor of researching Dark Matter. Unfortunately, his grant-funded position at Caltech is dependent on his continued work in String Theory, so he must search elsewhere. He applies and receives offers from MIT and Harvard. While money is no object to Sheldon, he wants to ensure he's paid fairly and that his offers are at least the median salary among the two schools' Physics departments. Therefore, he hires you to find the median salary across the two departments. Each school maintains a database of all of the salaries for that particular school, but there is no central database.

Each school has given you the ability to access their particular data by executing *queries*. For each query, you provide a particular database with a value k such that $1 \leq k \leq n$, and the database returns to you the k^{th} smallest salary in that school's Physics department.

You may assume that: each school has exactly n physicists (i.e. $2n$ total physicists across both schools), every salary is unique (i.e. no two physicists, regardless of school, have the same salary), and we define the *median* as the n^{th} highest salary across both schools.

1. Design an algorithm that finds the median salary across both schools in $\Theta(\log(n))$ total queries.
2. State the complete recurrence for your algorithm. You may put your $f(n)$ in big-theta notation. Show that the solution for your recurrence is $\Theta(\log(n))$.
3. Prove that your algorithm above finds the correct answer. *Hint: Do induction on the size of the input.*

Solution:

We have two databases, and we call MIT's database as data1 and Harvard as data2. We can search from the highest of data1 and lowest salary of data2. If the value of data1 is higher than data2, we query $n - \frac{n}{2}$ of data1, and $n + \frac{n}{2}$ of data2. We then recursively compare the values from two databases and make new query as (current position $-\frac{n}{2^i}$) with the data that this higher or (current position $+\frac{n}{2^i}$) to the data that is lower. Then we have the number closer and closer to the median until there is only one number of salary left for each database, and thus in the base case where $\frac{n}{2^i} = 1$, we return the average of that. It is the median.

```
def findMedian(posData1 = n, posData2 = 1, numofsalary = n)
    if numofsalary <= 1:
        query1 = query_data1(posData1)
        query2 = query_data2(posData2)
        return (query1+query2)/2
    query1 = query_data1(posData1)
    query2 = query_data2(posData2)
    if query1 > query2:
        if (posData1 - numofsalary/2) >= 1:
            return findMedian(posData1 - numofsalary/2, posData2 +
                               numofsalary/2, numofsalary/2)
        else:
            return (query1+query2)/2
    else:
        if (posData2 - numofsalary/2) >= 1:
            return findMedian(posData1 + numofsalary/2,
                               posData2 - numofsalary/2, numofsalary/2)
        else:
            return (query1+query2)/2
```


The recurrence has the following relation:

$$T(n) = T\left(\frac{n}{2}\right) + 2$$

On each recursive steps, we have a cost of 2 since we do one query that has a cost of 1 on two datasets. Thus, $f(n) \in O(\log n)$.

We have $\log n$ total levels of recursive steps, and on each level, we have one case. The total time is $T(n) = 2 \log n$. So, $T(n) \in \Theta(\log n)$.

Proof.

We prove this by induction. First, check the base case where the *numofsalary* $n = 1$. When $n = 1$, there are two salaries in total (one from each dataset). The function returns the average of the two salaries as the median, which is correct, and thus the base case is true.

Hence we assume that $\forall n < k$, $findMedian(n, 0, n)$ is correct. We now consider that the *numofsalary* $n = k$.

$findMedian(k, 0, k)$ returns $findMedian(k/2, k/2, k/2)$. As $n = k/2 < k$, $findMedian(k/2, k/2, k/2)$ is correct. Thus, $findMedian(k, 0, k)$ is correct.

As the base case is true and when the case with n is true, the case with $n + 1$ is also true. We can conclude by induction that the algorithm is true and returns the correct answer. \square

PROBLEM 6 *Mission Impossible*

As the newly-appointed Secretary of the Impossible Missions Force (IMF), you have been tasked with identifying the double agents that have infiltrated your ranks. There are currently n agents in your organization, and luckily, you know that the majority of them (i.e., strictly more than $n/2$ agents) are loyal to the IMF. All of your agents know who is loyal and who is a double agent. Your plan for identifying the double agents is to pair them up and ask each agent to identify whether the other is a double agent or not. Agents loyal to your organization will always answer honestly while double agents can answer arbitrarily. The list of potential responses are listed below:

Agent 001	Agent 002	Implication
"002 is a double agent"	"001 is a double agent"	At least one is a double agent
"002 is a double agent"	"001 is loyal"	At least one is a double agent
"002 is loyal"	"001 is a double agent"	At least one is a double agent
"002 is loyal"	"001 is loyal"	Both are loyal or both are double agents

1. A group of n agents is "acceptable" for a mission if a majority of them ($> n/2$) are loyal. Suppose we have an "acceptable" group of n agents. Describe an algorithm that has the following properties:
 - Uses at most $\lfloor n/2 \rfloor$ pairwise tests between agents.
 - Outputs a smaller "acceptable" group of agents of size at most $\lceil n/2 \rceil$.
2. Using your approach from Part 1, devise an algorithm that identifies which agents are loyal and which are double agents using $\Theta(n)$ pairwise tests.

Solution:

1:

We put all agents into two groups and paired them into groups of two (single agent left out and do nothing if the total number is odd). So we now have $\frac{n}{2}$ pairs of paired agents. Since there are more than half of agents are loyal agents, we do the following (assume a pair is (A, B)):

(i) If agent B says that agent A is loyal, we put agent A into the *loyalset1* no matter whether B is loyal or not. If agent A says B is double, then do nothing with B, no matter whether A is loyal or not.

(ii) If agent A says that agent B is loyal, we put agent B into the *loyalset2* no matter whether A is loyal or not. If agent B says A is double, then do nothing with A, no matter whether B is loyal or not.

Lastly, we get two loyal sets, and we choose the smaller one to do the mission. We add the left-out agent into the group if the total number is odd. We choose either group if the number of people in each group is the same. It works because we have more than half of loyal agents among all agents.

2:

We run the recursive steps from 1 repeatedly on each smaller group to eventually get one loyal agent. This final agent is highly likely to be loyal, and thus we can pair him with all the rest agents to use my loyal agents to spot all the double agents.

We have $n - 1$ pairs. The time complexity is: $\frac{n}{2} + \frac{n}{4} + \frac{n}{8} + \dots + 1 + (n - 1) = 2n - 1 \in \Theta(n)$.